

Principal components analysis and Factorial analysis to measure latent variables in a quantitative research: A mathematical theoretical approach

Arturo, GARCÍA-SANTILLÁN,
Milka, ESCALERA-CHÁVEZ,
Francisco, VENEGAS-MARTÍNEZ

ABSTRACT. The aim of this paper focuses on showing how the factorial analysis and principal components analysis are useful for measuring latent variables in a concise way and safely as a help to building for new concepts and theories.

1. PRELIMINARY NOTES AND NOTATION

In words of Wulder [5] “Multivariate statistics provide the ability to analyze a complex set of data. Principal components analysis (PCA) and factor analysis (FA) are statistical techniques applied to a single set of variables to discover which sets of variables in the set form coherent subsets that are relatively independent of one another. Variables that are correlated with one another which are also largely independent of other subsets of variables are combined into factors. Factors which are generated are thought to be representative of the underlying processes that have created the correlations among variables”.

Factor analysis is a multivariate statistical technique which allows obtaining a structure of latent variables in a data matrix known as factors therefore is considered as a data reduction technique conditionally if, his hypotheses are met; the information contained in the matrix can be expressed without a lot deviation, a lower number of dimensions represented by these factors [4].

Principal Component Analysis and Factor Analysis can be exploratory in nature; FA is used as a tool in attempts to reduce a large set of variables to a more meaningful, smaller set of variables. As both FA and PCA are sensitive to the magnitude of correlations robust comparisons must be made to ensure the quality of the analysis.

Correlation coefficients tend to be less reliable when estimated from small sample sizes. In general it is a minimum to have at least five cases for each observed variable. Missing data need be dealt with to provide for the best possible relationships between variables. Fitting missing data through regression techniques are likely to over fit the data and result in correlations to be unrealistically high and may as a result manufacture factors.

Normality provides for an enhanced solution, but some inference may still be derived from nonnormal data. Multivariate normality also implies that the relationships between variables are linear.

Univariate and multivariate outliers need to be screened out due to a heavy influence upon the calculation of correlation coefficients, which in turn has a strong influence on the calculation of factors. In PCA multicollinearity is not a problem as matrix inversion is not required, yet for most forms of FA singularity and multicollinearity is a problem.

If the determinant of R and eigenvalues associated with some factors approach zero, multicollinearity or singularity may be present. Deletion of singular or multicollinear variables is required.

The theorems and it implications are following: In order to measure; $X_1 X_2 \dots X_n$ observed random variables, which are defined in the same population that share, m ($m < p$) commons causes to find $m+p$ new variables, which we call common factors ($Z_1, Z_2, \dots Z_m$), besides, unique factors ($\varepsilon_1 \varepsilon_2 \dots \varepsilon_p$) in order to determine their contribution in the original variables ($X_1 X_2 \dots X_{p-1} X_p$) the model is now defined by the following equations according to Carrasco-Arroyo [1], Garcia-Santillán, Venegas-Martínez and Escalera-Chávez [3]:

$$\begin{array}{l}
 X_1 = a_{11}Z_1 + a_{12}Z_2 + \dots + a_{1m}Z_m + b_1\xi_1 \\
 X_2 = a_{21}Z_1 + a_{22}Z_2 + \dots + a_{2m}Z_m + b_2\xi_2 \\
 \dots \\
 X_p = a_{p1}Z_1 + a_{p2}Z_2 + \dots + a_{pm}Z_m + b_p\xi_p
 \end{array} \quad (1)$$

Where:

Z_1, Z_2, \dots, Z_m are common factors

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$ are unique factors

Thus, $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$ have influence in all variables X_i ($i=1, \dots, p$) ξ_i influence in X_i ($i=1, \dots, p$)

Model equations can be expressed in matrix form as follow:

$$\begin{array}{l}
 X_1 = \\
 X_2 = \\
 \dots \\
 X_p =
 \end{array}
 =
 \begin{array}{l}
 \left| \begin{array}{cccc}
 a_{11} & a_{12} & \dots & a_{1m} \\
 a_{21} & a_{22} & \dots & a_{2m} \\
 \dots & \dots & \dots & \dots \\
 a_{p1} & a_{p2} & \dots & a_{pm}
 \end{array} \right|
 \begin{array}{l}
 Z_1 \\
 Z_2 \\
 \dots \\
 Z_m
 \end{array}
 +
 \begin{array}{l}
 b_1 \xi_1 \\
 b_2 \xi_2 \\
 \dots \\
 b_p \xi_p
 \end{array}
 \end{array} \quad (2)$$

Therefore, the resulting model can be expressed in a condensed form as:

$$X = AZ + \xi_I \quad (3)$$

Where, we assume that $m < p$ because they want to explain the variables through a small number of new random variables and all of the $(m + p)$ factors are correlated variables, that is, that the variability explained by a variable factor, have not relation with the other factors.

We know that the each observed variable of model is a result of lineal combination of each common factor with different weights (a_{ia}), those weights are called saturations, but one of part of x_i is not explained for common factors.

As we know, all problems intuitive can be inconsistent when obtaining solutions and therefore, we require the approach of hypothesis; hence, in the factor model we used the following assumptions:

H₁: The factors are typified random variables, and inter correlated, like:

$$\begin{array}{l}
 E[Z_i] = 0 \quad E[\xi_i] = 0 \quad E[Z_i Z_i] = 1 \\
 E[\xi_i \xi_i] = 1 \quad E[Z_i Z_{i'}] = 0 \quad E[\xi_i \xi_{i'}] = 0 \\
 E[Z_i \xi_i] = 0
 \end{array}$$

Further, we must consider that the factors have a primary goal to study and simplify the correlations between variables, measures, through the correlation matrix, then, we will understand that:

H₂: The original variables could be typified by transforming these variables of type

$$X_i = \frac{X_i - \bar{X}}{\sigma_x}$$

Therefore, and considering the variance property we have:

$$\text{var}(x_i) = a_{i1}^2 \text{var}(z_1) + a_{i2}^2 \text{var}(z_2) + \dots + a_{im}^2 \text{var}(z_m) + b_i^2 \text{var}(\xi_i) \tag{5}$$

Resulting $1 = a_{i1}^2 + a_{i2}^2 + a_{i3}^2 + \dots + a_{im}^2 + b_i^2 \rightarrow \forall_i = 1 \dots p$ (6)

2. THEOREM SATURATIONS, COMMUNALITIES AND UNIQUENESS

We denominated *saturations* of the variable x_i in the factor z_a of coefficient a_{ia} In order to inform the relationship between the variables and the common factors is necessary determining the coefficient de A (assuming the hypotheses H_1 y H_2), where V is the matrix of eigenvectors and Λ matrix eigenvalues, so we obtained:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2m} \\ a_{31} & a_{32} & a_{33} & \dots & a_{3m} \\ \dots & \dots & \dots & \dots & \dots \\ a_{p1} & a_{p2} & a_{p3} & \dots & a_{pm} \end{bmatrix} \tag{7}$$

$$R = V\Lambda V^T = V\Lambda^{1/2}\Lambda^{1/2}V^T = AA^T, \tag{8}$$

$$A = V\Lambda^{1/2}$$

The above suggests that a_{ia} coincides with the correlation coefficient between the variables and factors. In the other sense, for the case of non-standardized variables, A is obtained from the covariance matrix S , hence the correlation between x_i and z_a is the ratio:

$$\text{corr}(i, a) = \frac{a_{ia}}{\sigma_a} = \frac{a_{ia}}{\sqrt{\lambda_a}} \tag{9}$$

Thus, the variance of the a_i factor is results of the sum of squares of saturations of a_i column of A :

$$\lambda_a = \sum_{i=1}^p a_{ia}^2 \tag{10}$$

Considering that:

$$A^T A = (V\Lambda^{1/2})^T (V\Lambda^{1/2}) = \Lambda^{1/2} V^T V \Lambda^{1/2} = \Lambda^{1/2} I \Lambda^{1/2} = \Lambda \tag{11}$$

We denominated communalities to the next theorem:

$$h_i^2 = \sum_{a=1}^m a_{ia}^2 \tag{12}$$

The communalities show a percentage of variance of each variable (i) that explains for m factors.

Thus, every coefficient h_i^2 is called variable specificity. Therefore the matrix model $\mathbf{X}=\mathbf{AZ}+\xi, \xi$ (unique factors matrix), \mathbf{Z} (common factors matrix) will be lower while greater be the variation explain for every m (common factor).

So, if we work with typified variables and considering the variance property, so, we have:

$$\begin{aligned} 1 &= a_{i1}^2 + a_{i2}^2 + \dots + a_{ia}^2 + b_i^2 \\ 1 &= h_i^2 + b_i^2 \end{aligned} \tag{13}$$

Recall that the variance of any variable, is the result of adding their communalities and the uniqueness b^2 , thus, the number of factors obtained, there is a part of the variability of the original variables unexplained and correspond to a residue (unique factor).

Reduced correlation matrix

Based on correlation between variables i and i' we have now:

$$\text{corr}(x_i, x_{i'}) = \frac{\text{cov}(x_i, x_{i'})}{\sigma_i \sigma_{i'}} \tag{14}$$

Also, we know

$$x_i = \sum_{a=1}^m a_{ia} z_a + b_i \varepsilon_i, \quad x_{i'} = \sum_{a=1}^m a_{i'a} z_a + b_{i'} \varepsilon_{i'} \tag{15}$$

The hypothesis which we started, now we have:

$$\text{corr}(x_i, x_{i'}) = \text{cov}(x_i, x_{i'}) = \sigma_{ii'} = E \left[\left(\sum_{a=1}^m a_{ia} z_a + b_i \varepsilon_i \right) \left(\sum_{a=1}^m a_{i'a} z_a + b_{i'} \varepsilon_{i'} \right) \right] \tag{16}$$

Developing the product:

$$= E \left[\sum_{a=1}^m a_{ia} a_{i'a} z_a z_a + \sum_{a=1}^m a_{ia} b_{i'} z_a \varepsilon_{i'} + \sum_{a=1}^m b_i a_{i'} \varepsilon_i z_a + \sum_{a=1}^m b_i b_{i'} \varepsilon_i \varepsilon_{i'} \right] \tag{17}$$

From the linearity of hope and considering that the factors are uncorrelated (hypotheses of starting), now we have:

$$\begin{aligned} \text{cov}(x_i, x_{i'}) &= \sigma_{ii'} = \sum_{a=1}^m a_{ia} a_{i'a} = \text{corr}(x_i, x_{i'}) \\ \forall i, i' &\rightarrow 1 \dots \dots \dots p \end{aligned} \tag{18}$$

The variance of variable i -^{esim} is given for:

$$\begin{aligned} \text{var}(x_i) &= \sigma_i^2 = E [x_i x_i] = 1 = E \left[\sum_{a=1}^m (a_{ia} z_a + b_i \varepsilon_i)^2 \right] = \\ &= E \left[\sum_{a=1}^m (a_{ia}^2 z_a^2 + b_i^2 \varepsilon_i^2 + 2a_{ia} b_i z_a \varepsilon_i) \right] \end{aligned} \tag{19}$$

If we take again the start hypothesis, we can prove the follow expression:

$$\sigma_i^2 = 1 = \sum_{a=1}^m a_{ia}^2 + b_i^2 = h_i^2 + b_i^2 \tag{20}$$

In this way, we can test how the variance is divided into two parts: the communality and uniqueness, which is the residual variance not explained by the model

Therefore, we can say that the matrix form is: $\mathbf{R} = \mathbf{A}\mathbf{A}' + \xi$ where $\mathbf{R}^* = \mathbf{R} - \xi^2$.

\mathbf{R}^* is a reproduced correlation matrix, obtained from the matrix \mathbf{R}

$$\mathbf{R}^* = \begin{bmatrix} h_1^2 & r_{12} & r_{13} & r_{14} & \dots & r_{1p} \\ r_{21} & h_2^2 & r_{23} & r_{24} & \dots & r_{2p} \\ r_{31} & r_{32} & h_3^2 & r_{34} & \dots & r_{3p} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ r_{p1} & r_{p2} & r_{p3} & r_{p4} & \dots & h_p^2 \end{bmatrix} \tag{21}$$

The fundamental identity is equivalent to the following expression: $\mathbf{R}^* = \mathbf{A}\mathbf{A}'$. Therefore the sample correlation matrix is a matrix estimator $\mathbf{A}\mathbf{A}'$. Meanwhile, a_{ia} saturation coefficients of variables in the factors, should verify this condition, which certainly, is not enough to determine them.

When the product is estimated $\mathbf{A}\mathbf{A}'$, we diagonalizable the reduced correlation matrix, whereas a solution of the equation would be: $\mathbf{R} - \xi^2 = \mathbf{R}^* = \mathbf{A}\mathbf{A}'$ is the matrix \mathbf{A} , whose columns are the standardized eigenvectors of \mathbf{R}^* . From this reduced matrix, through a diagonal, as mathematical instrument, we obtain through vectors and eigenvalues, the factor axes.

Factorial analysis viability

To validate the appropriateness of factorial model is necessary to design the sample correlation matrix \mathbf{R} , from the data obtained. Also be performed prior hypothesis tests to determine the relevance of the factor model, that is, whether it is appropriate to analyze the data with this model [2,3].

A contrast to be performed is the Bartlett Test of Sphericity. It seeks to determine whether there is a relationship structure –relationships-- or not among the original variables. The correlation matrix \mathbf{R} indicates the relationship between each pair of variables (r_{ij}) and its diagonal will be compose for 1(ones).

Hence, if there is not relationship between the variables h , then, all correlation coefficients between each pair of variable would be zero. Therefore, the population correlation matrix coincides with the identity matrix and determinant will be equal to 1.

$$\begin{aligned} H_0 : |\mathbf{R}| &= 1 \\ H_1 : |\mathbf{R}| &\neq 1 \end{aligned}$$

If the data are a random sample from a multivariate normal distribution, then, under the null hypothesis, the determinant of the matrix is 1 and is shown as follows:

$$-\left[n - 1 - \frac{(2p + 5)}{6} \right] \ln |\mathbf{R}| \tag{22}$$

Under the null hypothesis, this statistic is asymptotically distributed through a χ^2 distribution with $p(p-1)/2$ degrees freedom. So, in case of accepting the null hypothesis would not be advisable to perform factor analysis.

Therefore, the equation can be expressed as follows:

$$x = Af + u \hat{U} X = FA' + U \tag{25}$$

Where:

<i>Data matrix</i>	<i>Factorial load matrix</i>	<i>Factorial matrix</i>
$x = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_p \end{pmatrix}, f = \begin{pmatrix} F_1 \\ F_2 \\ \dots \\ F_k \end{pmatrix}, u = \begin{pmatrix} u_1 \\ u_2 \\ \dots \\ u_p \end{pmatrix}$	$A = \begin{pmatrix} a_{11} a_{12} \dots a_{1k} \\ a_{21} a_{22} \dots a_{2k} \\ \dots \\ a_{p1} a_{p2} \dots a_{pk} \end{pmatrix}$	$F = \begin{pmatrix} f_{11} f_{12} \dots f_{1k} \\ f_{21} f_{22} \dots f_{2k} \\ \dots \\ f_{p1} f_{p2} \dots f_{pk} \end{pmatrix}$

With a variance equal to:

$$\text{Var}(X_i) = \sum_{j=1}^k a_{ij}^2 + \Psi_i = h_i^2 + \Psi_i; i = 1, \dots, p \tag{26}$$

Where:

$$h_i^2 = \text{Var} \left(\sum_{j=1}^k a_{ij} F_j \right) \dots y \dots \Psi_i = \text{Var}(u_i) \tag{27}$$

This equation corresponds to the communalities and the specificity of the variable X_i . Thus the variance of each variable can be divided into two parts:

- a) in their communalities h_i^2 representing the variance explained by common factors, and
- b) the specificity Ψ_i that represents the specific variance of each variable.

Thus obtaining:

$$\text{Cov}(X_i, X_l) = \text{Cov} \left(\sum_{j=1}^k a_{ij} F_j, \sum_{j=1}^k a_{lj} F_j \right) = \sum_{j=1}^k a_{ij} a_{lj} \quad \forall i \neq l \tag{28}$$

With the transformation of the correlation matrix's determinants, we obtained Bartlett's test of sphericity, and it is given by the following equation:

$$d_R = - \left[n - 1 - \frac{1}{6} (2p + 5) \ln |R| \right] = - \left[n - \frac{2p + 11}{6} \right] \sum_{j=1}^p \log(\lambda_j) \tag{29}$$

$$\left[n - \frac{2p + 11}{6} \right] \log \frac{\left[\frac{1}{p - m} \left(\text{traz} R^* - \left(\sum_{a=1}^m \lambda_a \right) \right) \right]^{p-m}}{|R^*| \prod_{a=1}^m \lambda_a} \tag{30}$$

Finally, to calculate principal components: For all cases we will have “ p ” initials variables:

$$X' = [X_1, X_2, \dots, X_p] \tag{31}$$

Thus, we build p principal components guided by:

- (1) Linear function of the original variables,
- (2) Absorbing the maximum variation of the variables X and
- (3) That are uncorrelated.

$$Y_{ij} = \hat{\beta}_{i1}X_{1j} + \hat{\beta}_{i2}X_{2j} + \dots + \hat{\beta}_{ip}X_{pj} ; \quad j = 1, 2, \dots, n \tag{32}$$

$$Y_i = X \hat{\beta}_i \tag{33}$$

$$Y_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ \cdot \\ \cdot \\ Y_{in} \end{bmatrix} ; \quad X = \begin{bmatrix} X_{11} & X_{21} & \dots & X_{p1} \\ X_{12} & X_{22} & \dots & X_{p2} \\ \dots & \dots & \dots & \dots \\ X_{1n} & X_{2n} & \dots & X_{pn} \end{bmatrix} ; \quad \hat{\beta}_i = \begin{bmatrix} \hat{\beta}_{i1} \\ \hat{\beta}_{i2} \\ \cdot \\ \cdot \\ \hat{\beta}_{ip} \end{bmatrix} \tag{34}$$

The variation of variable Y_i , will be:

$$Y_i' Y_i = \hat{\beta}_i' S \hat{\beta}_i \tag{35}$$

Where: $S = X'X$

In order to obtain the first and the second component, we have the following procedure:

The first component is:

$$Y_1 = X \hat{\beta}_1 \tag{35.1}$$

So we need seek to maximize:

$$Y_1' Y_1 = \hat{\beta}_1' S \hat{\beta}_1 \tag{35.2}$$

And to address the process we must require:

$$\hat{\beta}_1' \hat{\beta}_1 = 1 \dots \tag{35.3}$$

Therefore, to the end:

$$\begin{aligned} \text{Max } Z &= \hat{\beta}_1' S \hat{\beta}_1 - \hat{\lambda}_1 (\hat{\beta}_1' \hat{\beta}_1 - 1) \dots \text{ie:} \\ \frac{\partial Z}{\partial \hat{\beta}_1} &= 2S \hat{\beta}_1 - 2 \hat{\lambda}_1 \hat{\beta}_1 = 0 \end{aligned} \tag{35.4}$$

$$\begin{aligned} S \hat{\beta}_1 - \hat{\lambda}_1 \hat{\beta}_1 &= 0 \\ (S - \hat{\lambda}_1 I) \hat{\beta}_1 &= 0 \end{aligned} \tag{35.5}$$

Leaving the trivial solution we have: $|S - \hat{\lambda}_1 I| = 0$ starting from here, we found $\hat{\lambda}_1$ that
 $(S - \hat{\lambda}_1 I)\hat{\beta}_1 = 0$ substituted at gives us $\hat{\beta}_1$

The second component is:

$$Y_2 = X \hat{\beta}_2 \tag{36}$$

And once again we need seek to maximize

$$Y_2' Y_2 = \hat{\beta}_2' S \hat{\beta}_2 \tag{36.1}$$

Once again subject to

$$\hat{\beta}_2' \hat{\beta}_2 = 1 \tag{36.2}$$

To which we now add the lack of correlation with the first component: $Y_2' Y_1 = 0 \dots$
 Which equal

$$\hat{\beta}_2' S \hat{\beta}_1 = 0 \tag{36.3}$$

Which may be written as well as

$$\hat{\beta}_2' \hat{\beta}_1 = 0 \tag{36.4}$$

Therefore, the function to maximize is:

$$\text{Max } Z = \hat{\beta}_2' S \hat{\beta}_2 - \hat{\lambda}_2 (\hat{\beta}_2' \hat{\beta}_2 - 1) - \mu_1 (\hat{\beta}_2' \hat{\beta}_1) \tag{37}$$

After finding the first derivative and carrying out a series of reductions, we have:

$$S \hat{\beta}_2 - \hat{\lambda}_2 \hat{\beta}_2 = 0 \tag{38}$$

$$\text{ie..... } (S - \hat{\lambda}_2 I) \hat{\beta}_2 = 0 \tag{38.1}$$

For all cases, the third component and subsequences is:

$$Y_i = X \hat{\beta}_i \tag{39}$$

And once again we need seek to maximize

$$Y_i' \dots Y_2 Y_1 = \hat{\beta}_i' S \hat{\beta}_i \tag{39.1}$$

Once again subject to

$$\hat{\beta}'_i \hat{\beta}_i = 1 \quad (39.2)$$

Which we now add the lack of correlation with the first, second and subsequent component:

$$Y_i \dots Y_2 Y_1 = 0 \dots$$

Which equal to...

$$\hat{\beta}'_i \hat{\beta}'_2 S \hat{\beta}_1 = 0 \quad (39.3)$$

Which may be written as well as

$$\hat{\beta}'_i \hat{\beta}'_2 \hat{\beta}_1 = 0 \quad (39.4)$$

3. CONCLUSION

This demonstration reveals the importance and convenience of having prior knowledge of the factors being measured to elect to measure the latent variables upon which to base new concepts or supporting established theories.

REFERENCES

- [1] Carrasco-Arroyo, S. (2012) Apuntes de Análisis Factorial [Notes for Analysis Factorial] Universidad de Valencia.
- [2] García-Santillán, A. et al (2013). Cognitive, Affective and Behavioral Components that explain Attitude toward Statistics. Journal of Mathematics Research. Vol 4, No. 5, October 2012 pp 8-16. (2012) URL:<http://dx.doi.org/10.5539/jmr.v4n5p8>
- [3] García-Santillán, Venegas-Martínez and Escalera-Chávez (2013). An exploratory factorial analysis to measure attitude toward statistic. Empirical study in undergraduate students in a private university. International Journal of Research and Reviews in Applied Sciences, 14, (2), 356-366.
- [4] Hair, J. Anderson, R. (2001). Análisis Multivariante. Prentice Hall
- [5] Wulder, M. (2013) A Practical Guide to the Use of Selected Multivariate Statistics. Canadian Forest Service, Pacific Forestry Centre, Victoria, Canada.